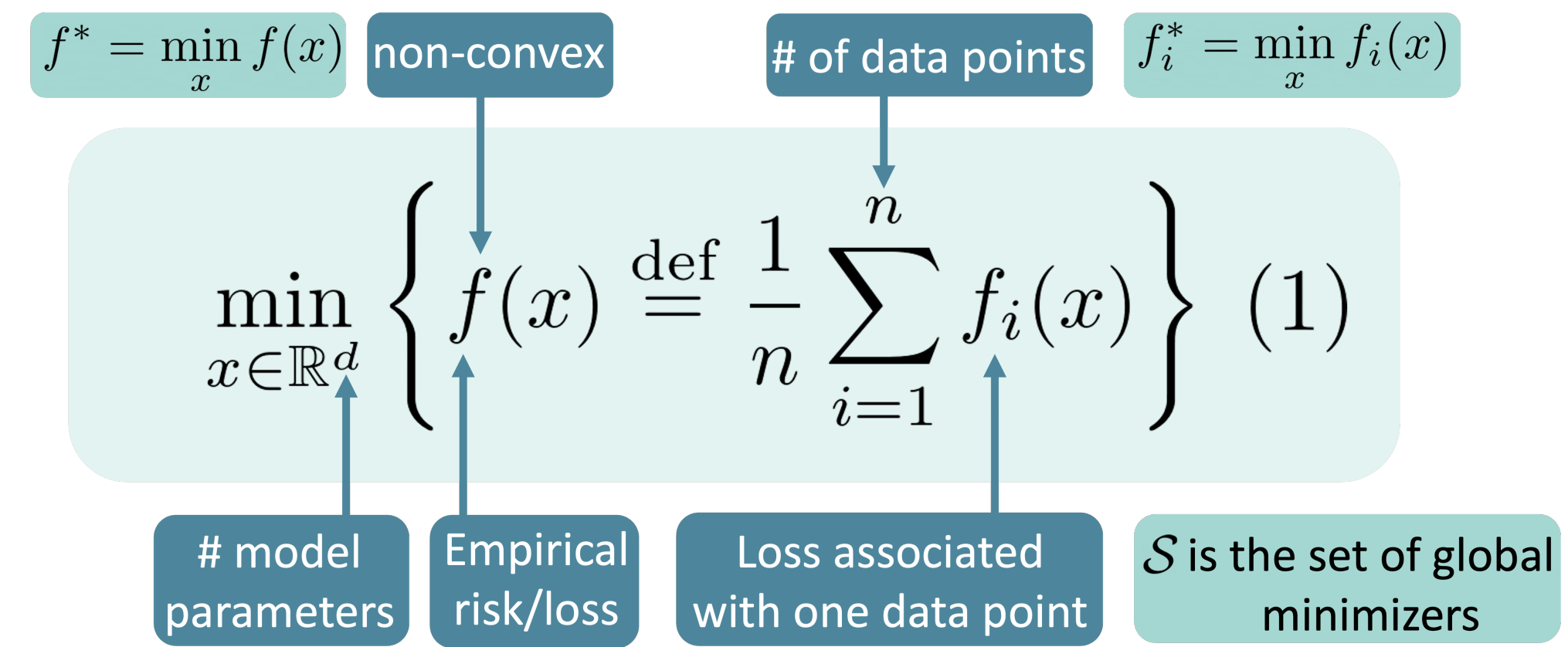




Problem Formulation

We want to solve a finite-sum problem



- To solve (1), optimizer's hyperparameters should be tuned to deliver a good performance in practice.
- Tuning of optimizers for training large models is an expensive and a time-consuming process.

These issues motivate us to design algorithms that are robust to the choice of hyperparameters, such as the learning rate.

NGN Step-size and Momentum

How to Add Coordinate-wise Adaptivity? For positive objective $f: \mathbb{R}^d \rightarrow \mathbb{R}$ we can approximate it around the point x as $f(x) = r^2(x) \Rightarrow f(x + \Delta) = r^2(x + \Delta) \approx (r(x) + \langle \nabla r(x), \Delta \rangle)^2$.

This leads to NGN algorithm design

$$x^{k+1} = \arg \min_{x \in \mathbb{R}^d} (r(x^k) + \langle \nabla r(x^k), x - x^k \rangle)^2 + \frac{1}{2\gamma} \|x - x^k\|^2 = x^k - \frac{\gamma}{1 + \frac{\gamma}{2f(x^k)} \|\nabla f(x^k)\|^2} \nabla f(x^k).$$

NGN converges for any stepsize hyperparameter γ for convex and smooth objectives $\{f_i\}_{i \in [n]}$.

How to Add Momentum? We consider the two most popular versions of incorporating momentum to NGN algorithm

$$\text{Version 1: } \begin{cases} \gamma_k = \frac{\gamma}{1 + \frac{\gamma}{2f_{S_k}(x^k)} \|\nabla f_{S_k}(x^k)\|^2} \\ m^k = \beta m^{k-1} + (1 - \beta) \gamma_k \nabla f_{S_k}(x^k) \quad (\text{Heavy-ball}) \\ x^{k+1} = x^k - m^k \end{cases}$$

$$\text{Version 2: } \begin{cases} \gamma_k = \frac{\gamma}{1 + \frac{\gamma}{2f_{S_k}(x^k)} \|m^k\|^2} \quad (\text{Adam-like}) \\ x^{k+1} = x^k - \gamma_k m^k \end{cases}$$

We aim for a version that preserves the robustness of NGN and achieves accelerated convergence on quadratic functions.

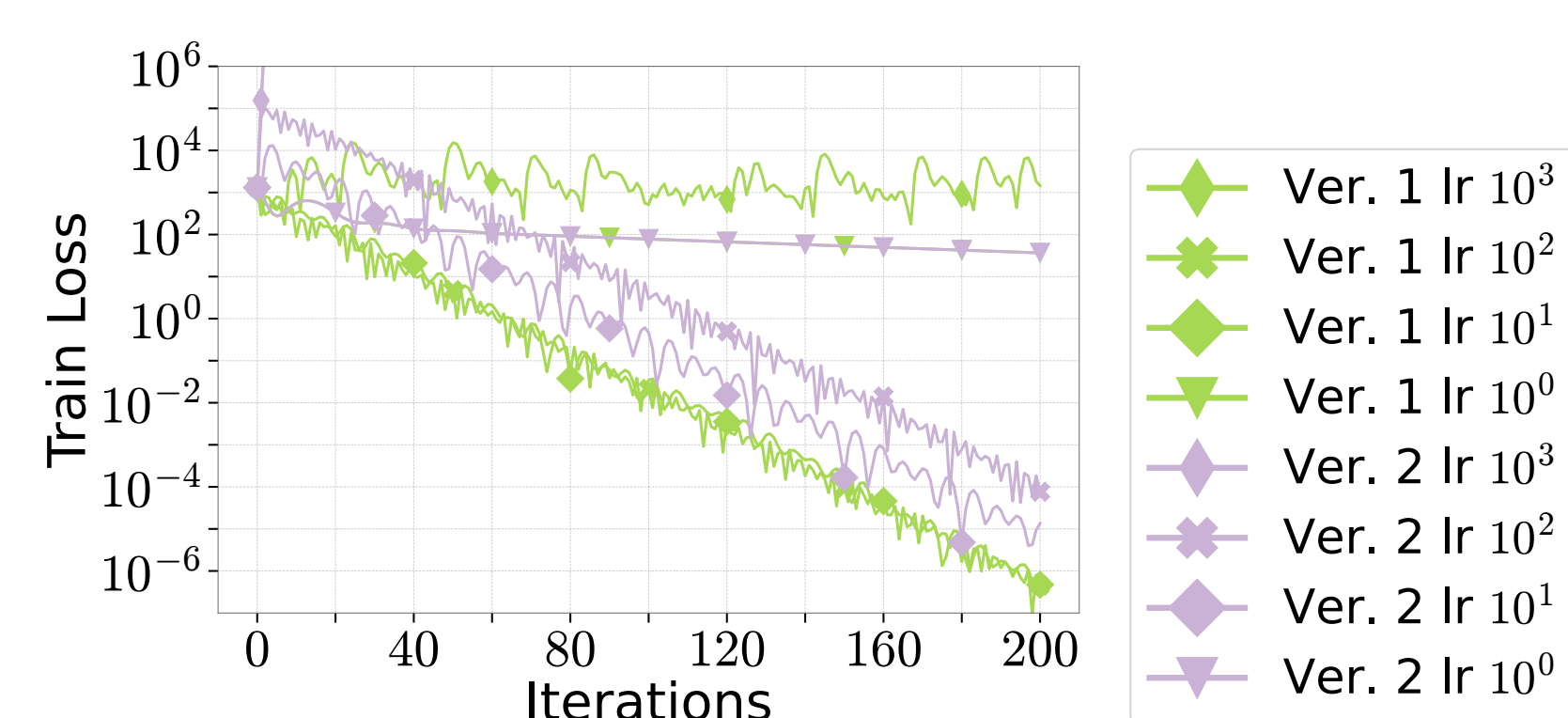


Figure 1: Comparison of two possible couplings of NGN and momentum.

Table 1: Summary of existing methods exploiting Polyak-type adaptive step-sizes and their convergence guarantees. **Mom.**—Supports momentum; **Diag.**—Supports diagonal step-sizes. $\sigma_{\text{int}}^2 := \mathbb{E}_S[f^* - f_S^*]$. \mathcal{O} notation hides absolute, problem-dependent constants and logarithmic factors.

Method	Rate ^(a)	Mom.	Diag.	Comments
SPS _{max} [Loizou et al., 2021]	$\mathcal{O}(1/K + \sigma_{\text{int}}^2)$	✗	✗	Conv. to non-vanishing neighbourhood
ALR-SMAG [Wang et al., 2023]	$\mathcal{O}((1 - \rho)^K + \sigma_{\text{int}}^2)$	✓	✗	Strong convexity Conv. to non-vanishing neighbourhood
Momo [Schaipp et al., 2024]	$\mathcal{O}(1/\sqrt{K})$	✓	✗	Bounded stoch. gradients Interpolation
Momo-Adam [Schaipp et al., 2024]	✗	✓	✓	Momo framework for Adam
MomSPS _{max} [Oikonomou and Loizou, 2024]	$\mathcal{O}(1/K + \sigma_{\text{int}}^2)$	✓	✗	Conv. to non-vanishing neighbourhood
NGN [Orvieto and Xiao, 2024]	$\mathcal{O}(1/\sqrt{K})$	✗	✗	—
IAM [Gower et al., 2025]	$\mathcal{O}(1/\sqrt{K})$	✓	✗	Knowledge of $f_i(x^*)$
NGN-M [This work]	$\mathcal{O}(1/\sqrt{K})$	✓	✗	—
NGN-MDv1 and NGN-MDv2 [This work]	✗	✓	✓	Combination of NGN-M and NGN-D
NGN-D [This work]	$\mathcal{O}(1/\sqrt{K})$	✗	✓	—

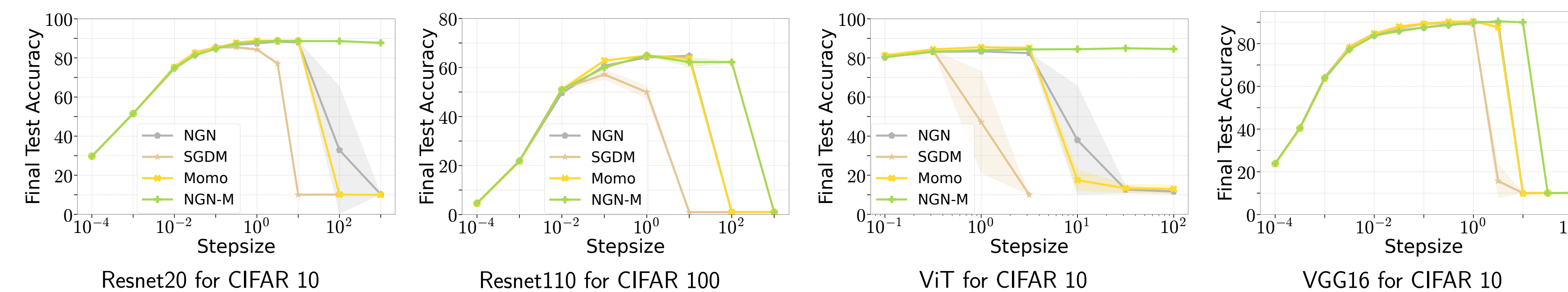


Figure 2: Comparison of NGN-M against several baselines with momentum when training small-scale models.

Theoretical Analysis of NGN-M

Convergence of NGN-M

Let f_i be positive, convex and L -smooth for all $i \in [n]$, and the errors $\sigma_{\text{int}}^2 := \mathbb{E}_S[f^* - f_S^*]$ $\sigma_{\text{pos}}^2 := \mathbb{E}_S[f_S^*]$ be bounded. Then we have

$$\mathbb{E}[f(\bar{x}^K) - f^*] \leq \frac{R(1 + 2\gamma L)^2}{\gamma K} + 8\gamma L(1 + 2\gamma L)^2 \sigma_{\text{int}}^2 + 2\gamma L \max\{0, 2\gamma L - 1\} \sigma_{\text{pos}}^2,$$

where $R := \|x^0 - x^*\|^2$ and $\bar{x}^K \sim \text{Unif}\{x^0, \dots, x^{K-1}\}$, and $\beta = \frac{\lambda}{1 + \lambda}$ with $\lambda \leq \min\{\gamma L, 0.5(1 + \gamma L)^{-1}(1 + 2\gamma L)^{-1}\}$.

- If we set $\gamma = \mathcal{O}(1/\sqrt{K})$, we obtain the convergence rate of order $\mathcal{O}(1/\sqrt{K})$.
- No restriction on the LR hyperparameter γ as long as the momentum parameter is small enough.
- In contrast to prior work, the convergence is obtained without strong assumptions such as interpolation condition $\sigma_{\text{int}}^2 = 0$.
- Small- β requirement can be moved to the LR hyperparameter γ under interpolation condition.
- Decaying LR hyperparameter γ as $\mathcal{O}(1/\sqrt{K})$ removes the need for knowing the total number of iterations K .

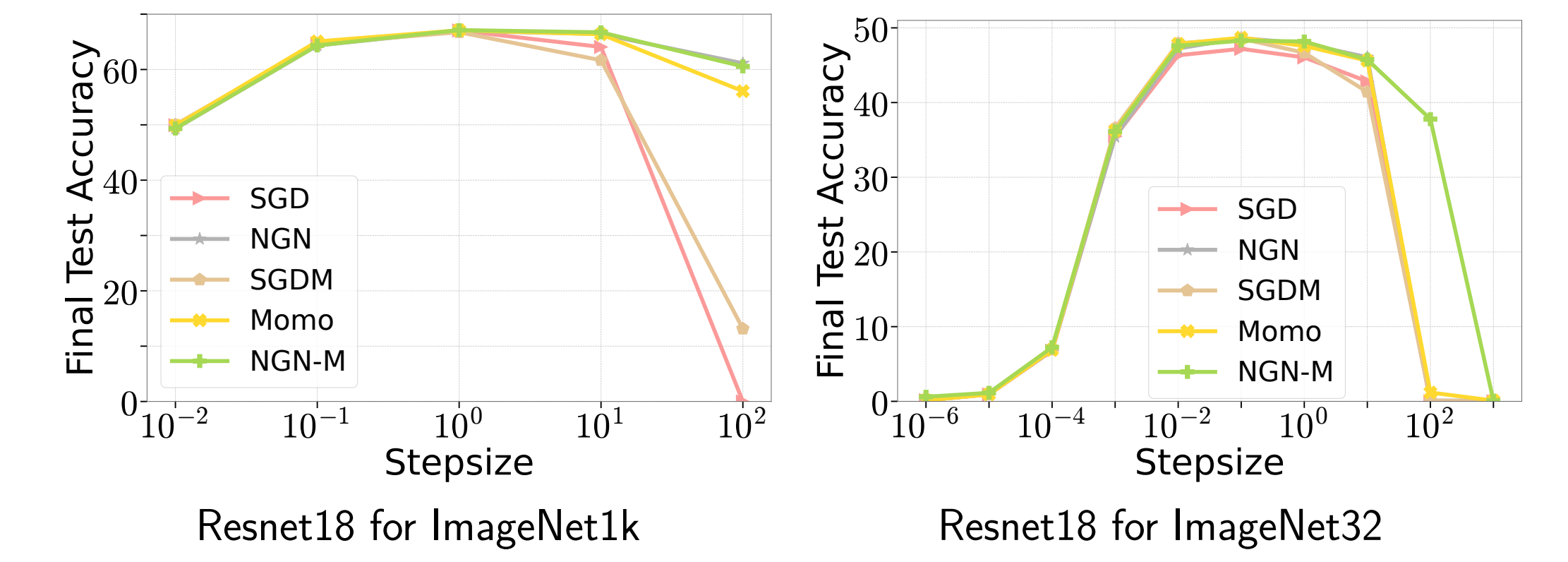


Figure 3: Comparison of NGN-M against several baselines with momentum when training Resnet18 on ImageNet datasets.

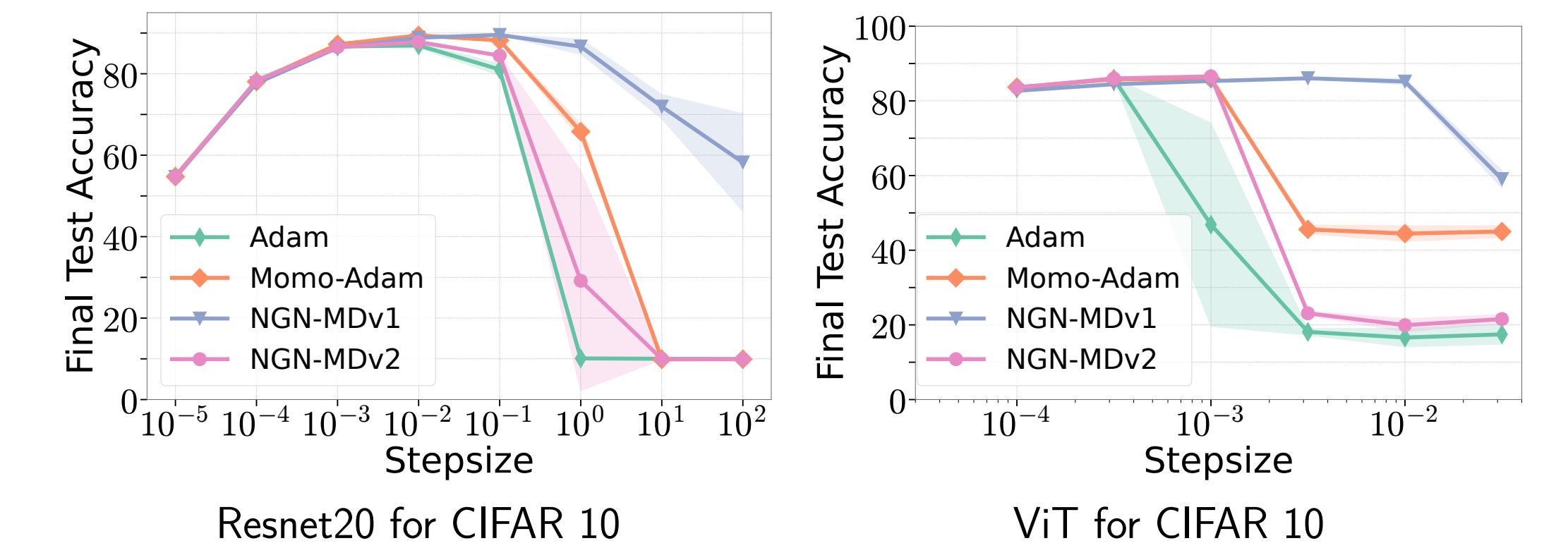


Figure 4: Comparison of NGN-MDv1 against several baselines with momentum and preconditioning when training Resnet18 on ImageNet datasets.

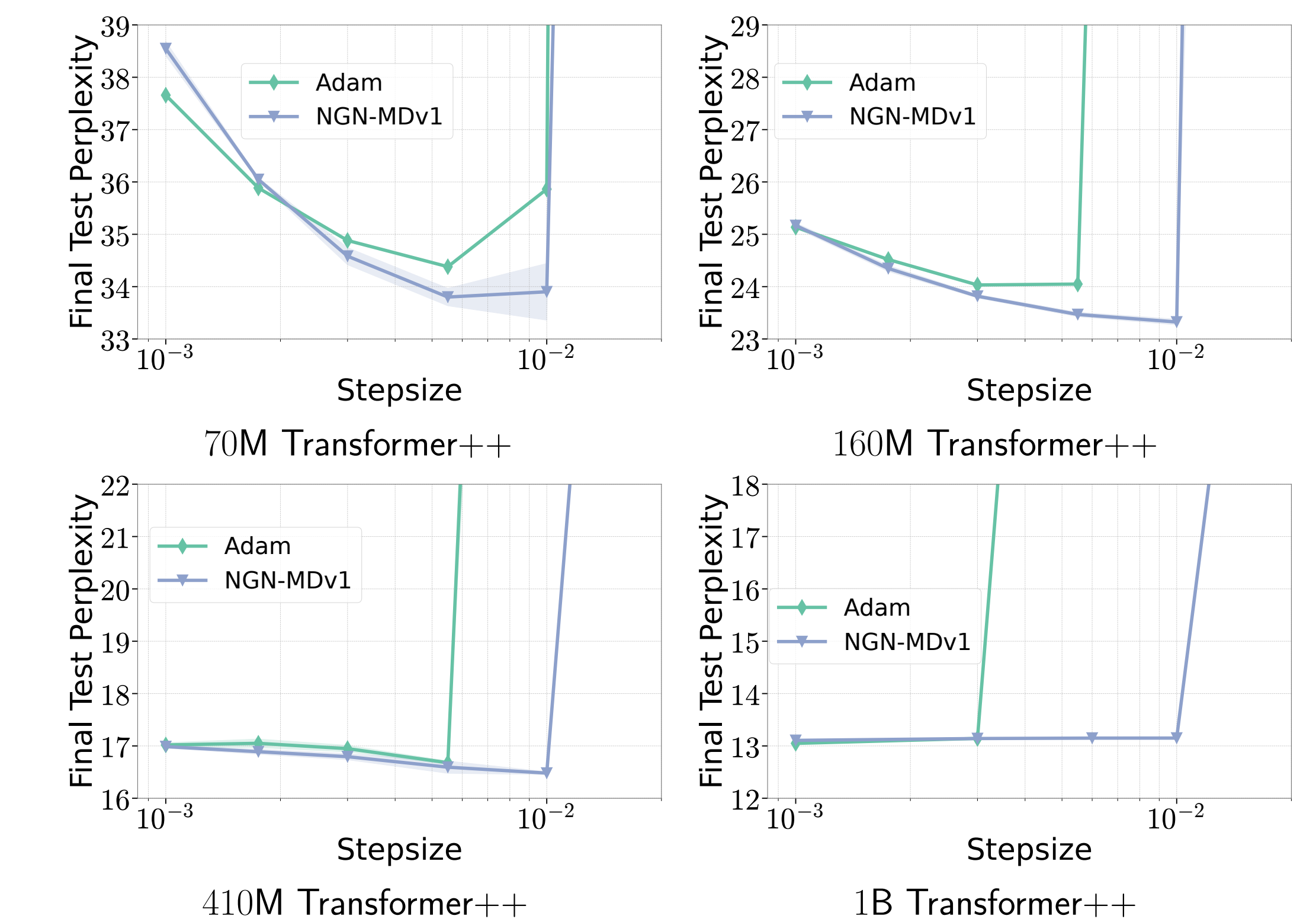


Figure 5: Comparison of stability to the LR hyperparameter across model sizes and optimizers in language modeling.

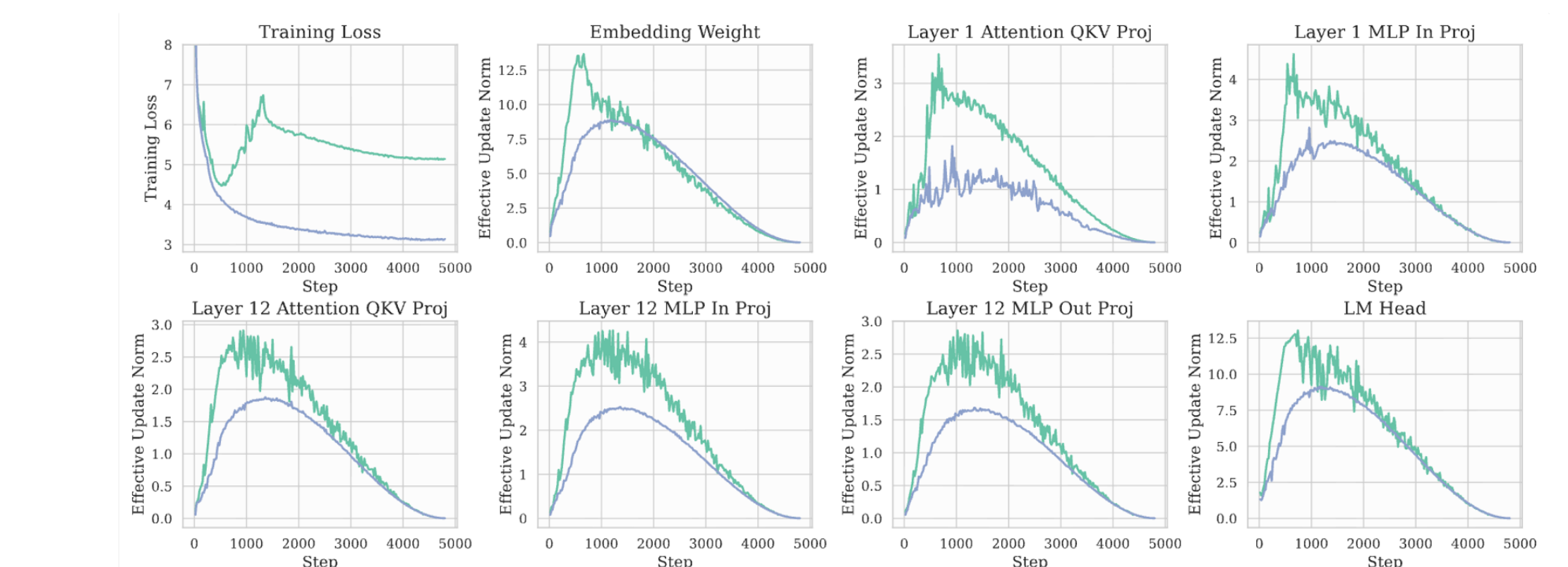


Figure 6: Magnitude of updates when training 160M language model with Adam and NGN-MDv1 and step-size hyperparameter 0.01.

Table 2: Train time of AdamW and NGN-MDv1 when training language models.

Model	Method	Time per Iteration (sec)	Time per Optimizer Update (sec)
70M	AdamW	1.63±0.01	0.0048±0.0002
	NGN-MDv1	1.65±0.01	0.0130±0.0002
160M	AdamW	3.33±0.03	0.0088±0.0003
	NGN-MDv1	3.37±0.02	0.0239±0.0003
410M	AdamW	8.41±0.06	0.0838±0.0009
	NGN-MDv1	8.68±0.06	0.2154±0.0007